# LLM Agent Meets Agentic AI: Can LLM Agents Simulate Customers to Evaluate Agentic-AI-based Shopping Assistants?

Lu Sun
University of California San Diego
La Jolla, California, United States
l5sun@ucsd.edu

Shihan Fu
Northeastern University
Boston, Massachusetts, United States
fu.shiha@northeastern.edu

Bingsheng Yao
Northeastern University
Boston, Massachusetts, United States
b.yao@northeastern.edu

Yuxuan Lu
Northeastern University
Boston, Massachusetts, United States
lu.yuxuan@northeastern.edu

Wenbo Li
North Carolina State University
Raleigh, North Carolina, United States
wli55@ncsu.edu

Hansu Gu
Independent Researcher
Seattle, Washington, United States

Jiri Gesi
Independent Researcher
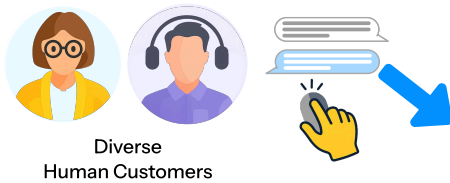Palo Alto, California, United States

Jing Huang
Independent Researcher
Palo Alto, California, United States

Chen Luo
Independent Researcher
Palo Alto, California, United States

Dakuo Wang*
Northeastern University
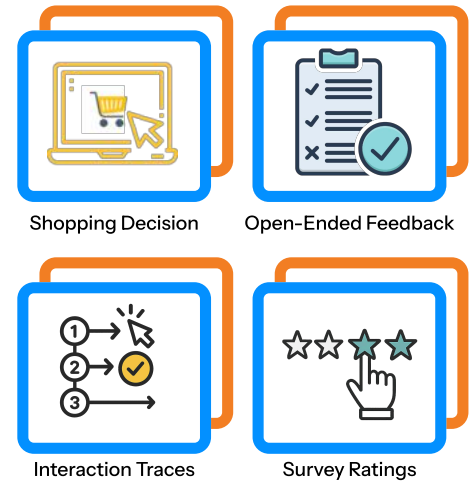Boston, Massachusetts, United States
d.wang@northeastern.edu

Figure 1: LLM Agent Meets Agentic AI: (a) We first conduct a human study to establish ground truth about how diverse human users interact with a conversational shopping assistant. (b) Building on these findings, we then simulate agentic AI agents as digital twins, role-playing the personas of real human participants to reproduce their decision-making behaviors in a controlled browser environment. (c) Finally, we evaluate human–agent alignment by comparing four aspects—shopping decisions, subjective feedback, interaction traces, and survey ratings—between real users and their simulated counterparts.

*Corresponding author

## Abstract

**Agentic AI** is emerging, capable of executing tasks through natural language, such as Copilot for coding or Amazon Rufus for shopping. Evaluating these systems is challenging, as their rapid evolution outpaces traditional human evaluation. Researchers have proposed **LLM Agents** to simulate participants as digital twins, but it remains unclear to what extent a digital twin can represent a specific customer in multi-turn interaction with an agentic AI system. In this paper, we recruited 40 human participants to shop with Amazon Rufus, collected their personas, interaction traces, and UX feedback, and then created digital twins to repeat the task. Pairwise comparison of human and digital-twin traces shows that while agents often explored more diverse choices, their action patterns aligned with humans and yielded similar design feedback. This study is the first to quantify how closely LLM agents can mirror human multi-turn interaction with an agentic AI system, highlighting their potential for scalable evaluation. Our code is open-sourced. [1]

## CCS Concepts

• **Do Not Use This Code** → **Generate the Correct Terms for Your Paper**.

## Keywords

LLM Agents

## 1 Introduction

Agentic AI systems are rapidly emerging across domains, enabling people to accomplish complex tasks through natural language interaction [4, 18]. From GitHub Copilot assisting developers in coding to OpenAI's ChatGPT supporting knowledge work, these systems represent a shift from static tools to proactive and conversational collaborators. A particularly impactful application area is e-commerce, where agentic AI transforms traditional search-and-filter shopping into dynamic multi-turn dialogues [49]. Instead of typing keywords into a search bar, customers can now interact with conversational shopping assistants that adapt to evolving preferences, refine recommendations in real time, and provide personalized guidance throughout the shopping journey. This agentic shift is exemplified by systems like Amazon's Rufus [49], Google Shopping's AI mode, and ChatGPT with integrated shopping capabilities [64], all of which demonstrate how AI-powered dialogue can reshape the online shopping experience into one that is more personalized, interactive, and engaging [18, 65].

Despite the rapid deployment and continual iteration of agentic-AI assistants, evaluating their performance and user experience (UX) has not kept pace. New AI features and capabilities can be released on a weekly or even daily basis, yet traditional evaluation methods such as think-aloud, heuristic walkthroughs, A/B testing,

or Wizard-of-Oz studies require weeks of planning, recruitment, and analysis [60, 79]. These methods have long provided valuable insights [38], but their slow and resource-intensive nature makes them poorly suited for assessing fast-evolving agentic AI systems in the real world [45]. The problem is compounded by the open-ended and adaptive nature of human-agentic-AI interactions: different users pursue different goals, and the system dynamically tailors its responses, introducing high variability that resists standardized evaluation criteria [47]. Taken together, these factors create a widening gap: while agentic-AI systems evolve rapidly, human-centered evaluation struggles to keep up, leaving researchers without scalable tools to capture how these systems shape decision-making, trust, and usability in practice.

To overcome the limitations of traditional UX evaluation methods, recent work has turned to large language model (LLM) agents as scalable evaluators [98, 102]. The prevailing **Agent-as-a-Judge** paradigm tasks LLMs with assessing the outputs of other AI models—for example, rating the quality of generated code or verifying factual accuracy, where evaluation criteria are relatively objective and framed in single-turn interactions [10, 102]. Similarly, researchers in the social sciences have used LLM agents to simulate survey respondents or role-play individuals, yielding plausible outcomes such as political attitudes or consumer preferences, yet still confined to single-turn responses. While these approaches demonstrate the potential of LLM agents as human **digital twins** for scalable evaluation, they remain restricted to judging static outputs or isolated prompts. Moreover, prior work has largely emphasized algorithmic correctness and agent performance [44, 58], overlooking user-centered dimensions such as trust, usability, and decision support that are critical for evaluating real-world interactions with agentic AI systems. What remains missing is the ability to evaluate **multi-turn human–AI interactions**, where users dynamically shape the trajectory of a conversation and the system adapts in real time. This gap motivates our central question: **can LLM agents go beyond judging isolated responses to role-play customers in dynamic multi-turn interaction with agentic AI systems?**

Recent studies have explored LLM-agent simulations in benchmark tasks [58, 101]. For example, the Mind2Web2 benchmark introduced 130 web tasks and constructed task-specific judge agents to automatically assess correctness on search tasks. These efforts highlight the potential of agent-based evaluation, but they generally emphasize final outcomes such as negotiated prices or task accuracy, and often operate in constrained environments disconnected from real users' multi-turn behavior. In this work, we narrow the focus to online shopping, a domain that is preference-rich, decision-intensive, and inherently conversational in real-world shopping. Our study centers on conversational shopping assistants, which can guide customers through discovery, comparison, and decision-making via dialogue. A leading example is Amazon Rufus [49], an already-deployed, commercially available system that illustrates how agentic AI can mediate the end-to-end shopping journey via conversational interaction. This setting is especially valuable for studying digital twins. It combines wide deployment with diverse users and supports multi-turn interactions that highlight conversation quality. It also produces observable outcomes such as product choices, dialogue trajectories, and UX ratings, which enable user-centered evaluation beyond algorithmic correctness.

To investigate this research question, we designed a two-stage evaluation pipeline that head-to-head compares a human customer's data with their LLM-agent digital twin's data. In the first stage, we conducted a large-scale user study in which 40 participants interacted with Amazon Rufus [2] to complete two representative shopping tasks. This produced a rich dataset of multi-turn conversations, product choices, and UX feedback. In the second stage, we instantiated persona-aligned LLM agents via UXAgent [45] to role-play as digital twins of the same participants and repeat the tasks. This approach enables pairwise comparisons of humans and their digital twins across decision outcomes, interaction behaviors, and evaluation results, providing a foundation for assessing how closely LLM agents can mirror real users in multi-turn interaction with agentic AI systems.

Our analyses reveal that LLM agents can meaningfully approximate human behavior while also exposing important gaps. First, agents consistently completed the shopping tasks, matching humans in overall interaction turn counts and final buy-or-not decisions (F1 score of 0.9). For example, the agent's opening queries showed some alignment with humans' first queries (similarity > 0.4). This result demonstrates that digital twins can capture the broad structure of human–AI interaction. Second, despite these similarities, trajectories soon diverged: sequence-level comparisons showed low overlap (similarity < 0.2), and only about 2% of agent–human pairs chose the same product. These differences highlight opportunities to improve how LLMs model human exploration and decision strategies via various fine-tuning techniques. Finally, when acting as evaluators, agents produced UX ratings aligned with human judgments on objective dimensions such as query relevance and coherence. Yet they tended to rate their own satisfaction more conservatively, while at the same time expressing a stronger preference for interacting with conversational shopping assistants over traditional search. Overall, these findings indicate that digital twins can reproduce many functional aspects of shopping interactions while also pointing to concrete directions for training them to better capture the nuance of human decision-making.

Building on these findings, our study demonstrates both the promise and the limitations of using persona-grounded LLM agents as digital twins for evaluating agentic AI systems. Although our analysis is situated in online shopping, the framework and insights extend to other applications where agentic AI engages users in adaptive, multi-turn interactions. More broadly, this work contributes to the methodological toolkit of human–AI interaction by combining the breadth of agent-based simulation with the depth of human-centered evaluation.

In sum, our contributions are:

- The first large-scale human study of multi-turn interactions with an agentic-AI–powered conversational shopping assistant (Amazon Rufus), capturing buy-or-not outcomes, interaction traces, and UX ratings.
- The first simulation framework that instantiates persona-grounded LLM agents as digital twins, enabling direct pairwise comparison with real customers across tasks, behaviors, and evaluations.

- Empirical insights into where agents align with or diverge from humans, providing the first evidence of both their potential as scalable evaluators and their limitations in capturing human-like reasoning and experience.

## 2 Related Work

### 2.1 From Traditional Online Shopping to Agentic-AI-based Shopping Assistants

Agentic AI systems are rapidly emerging and reshaping how users engage in online activities across multiple domains [1, 52, 67, 73]. For example, GitHub Copilot [19] is transforming software development through conversational code generation [86, 94], while Google Gemini [21] is redefining search by enabling natural language–based exploration [50, 74, 83]. These systems are capable of executing complex tasks through iterative interactions in natural language [78, 92], lowering barriers to access and broadening participation [70, 93].

Among these, Conversational Shopping Assistants (CSAs) have quickly become one of the most visible applications. In early 2025, several prominent companies—including Amazon [49], Google [48, 65], and OpenAI [64]—introduced their own versions of CSAs, which rapidly gained widespread attention and adoption. Traditionally, online shopping required users to rely on search boxes and filters, manually navigating across multiple tabs to browse, compare, and ultimately select products [27, 28, 51, 66]. By contrast, CSAs have begun to transform this process [4, 18]. Instead of keyword searches and static filters, users can now engage in multi-turn, natural language conversations with LLM-based CSAs to explore product options, refine their preferences, and receive tailored recommendations [91].



**Figure 2: Amazon Rufus Conversational Shopping Assistant. The user can send in a natural language query to Rufus. Rufus can respond with text responses, product recommendation cards, and related question suggestions.**

A typical Conversational Shopping Assistant (CSA) system like Rufus (Fig. 2) is seamlessly integrated into the shopping website interface Users can open the CSA through a pop-up window available from the navigation bar on any page, making the assistant easily accessible without interrupting the normal shopping flow. During browsing or product search, users can engage with the CSA at any time to request personalized product recommendations, seek clarifications, or compare alternatives. CSA supports natural language

---

queries (e.g., "I want to buy a monitor") and responds in a multi-modal format that includes explanatory text responses, product recommendation cards with images and pricing, as well as related question suggestions that guide further exploration. This interaction loop allows users to iteratively refine their preferences, verify product attributes, and make more informed purchase decisions, thereby augmenting the traditional search-and-filter paradigm with a conversational and context-aware shopping experience.

## 2.2 Challenges in Evaluating Agentic AIs

Evaluating the usability of agentic AI systems, especially those involving conversational assistants, presents unique methodological challenges [7, 11, 46]. First, agentic AI can exhibit highly adaptive behaviors [87]. Even when faced with the same task prompt, different users may choose distinct strategies, phrasing, or levels of detail in their inputs [85]. As a result, the same task can unfold differently across users and even across sessions with the same user [93]. This variability creates challenges for evaluation, as standardized usability measures—such as task completion time or error rates—may not capture the full spectrum of user behaviors and contexts [34, 99]. For example, two users might reach the same outcome through very different interaction paths, making it difficult to determine whether one experience is objectively "better" than the other. Second, users' perceptions of response quality and task support remain inherently subjective [15, 62]. These perceptions are shaped not only by the correctness or completeness of system outputs, but also by individual goals, prior knowledge, cognitive styles, and situational needs [96]. What one user considers a helpful explanation, another may view as excessive or confusing. Moreover, perceptions can shift dynamically during interaction: a user who initially values speed may later prioritize detail and transparency once they gain confidence in the system [35, 89]. Such subjective and evolving evaluations further complicate the assessment of system effectiveness, as measures of satisfaction, trust, or usefulness may vary widely both across and within users [32].

Traditional usability testing approaches often rely on human-centered methods that directly involve researchers, experts, or participants [23]. These methods are designed to uncover usability issues, evaluate user experience, and ensure that systems align with human needs [26, 68]. For example, expert heuristics [38, 55] involve evaluators applying established principles to identify usability flaws early in the design process, but their insights can be limited by the evaluators' expertise and may overlook domain-specific challenges [2, 3]. Controlled user studies [39, 72] provide empirical data through carefully designed tasks and measures, yet they require significant time, participant recruitment, and researcher effort to conduct, making them costly and difficult to repeat at scale [33]. Wizard-of-Oz methods [12, 37] allow researchers to prototype interactive systems by simulating system responses through human operators, offering rich insights into user expectations and interaction patterns, but they are labor-intensive and difficult to sustain beyond small studies [14, 31].

While these methods have been foundational in HCI and continue to offer valuable insights, they are resource-intensive and may not scale well for rapidly evolving agentic AI systems. This limitation is particularly evident in the context of LLM-based conversational assistants, where system behavior is highly dynamic, responses are non-deterministic, and user interactions can vary significantly across sessions [9, 36]. As a result, traditional approaches struggle to keep pace with the scale, speed, and variability of modern AI systems, highlighting the need for complementary methods of evaluation.

## 2.3 LLM Agents as Proxies for Human Participants in UX Evaluation

Since involving human participants in the evaluation of AI-powered systems requires substantial time and effort, researchers have increasingly explored using agents as evaluators of other systems. LLM-as-a-Judge [98] is one such approach, where an LLM is used as a judge to compare the outputs of different models on the same task. In benchmarks like Chatbot Arena [98] and MT-Bench [97], this method is commonly used to determine which response is better while significantly reducing the cost of human evaluation. Building on this idea, the Agent-as-a-Judge framework [102] extends LLM-as-a-Judge by equipping the judge with agentic capabilities such as autonomous planning, information retrieval, step-by-step execution, and tracking intermediate artifacts. This enables the judge to function as an agent and evaluate the entire task-solving process rather than just the final output, demonstrating that prompt-driven LLM agents can perform complex evaluation tasks traditionally handled by human experts. Meanwhile, broader evaluation suites such as AgentBench[42], MLAgentBench[30], and Mind2Web2[22] provide standardized tasks for measuring LLM agents' capabilities. For example, Mind2Web2[22] evaluates how well AI agents handle realistic, long-horizon, and dynamic web search tasks. It adopts an Agent-as-a-Judge framework, where a judge agent verifies whether each answer satisfies all task requirements and includes reliable citations. This setup allows researchers to systematically evaluate agentic search systems on complex real-world tasks. More specifically, in the context of online shopping, Zhu et al. [101] proposed an Agent-to-Agent simulation framework to systematically study how LLM-based agents negotiate and make decisions on behalf of human consumers or merchants. In this framework, LLM agents simulate buyers and sellers in consumer markets, engaging in multi-turn dialogues to negotiate prices and complete transactions. This allows the researchers to evaluate each agent's decision-making strategies and negotiation capabilities. However, such evaluations mainly focus on the final negotiation outcomes (e.g., agreed prices or profit margins) rather than the fine-grained interaction dynamics that typically occur when users interact with conversational shopping assistants. Yet, this line of research centers on technical performance, paying little attention to the diverse needs and intentions that shape how different user groups engage with AI-powered conversational systems.

In response, researchers have begun examining how LLM agents can generate actions that resemble those of diverse human users. [80, 90]. For example, LLM agents have been used to emulate a community of 25 residents in a virtual village [57], to replicate participants in social science studies [24, 40, 59, 69], to act as patients and clinicians in hospital contexts [41], and to take on the role of software developers in company settings [61].
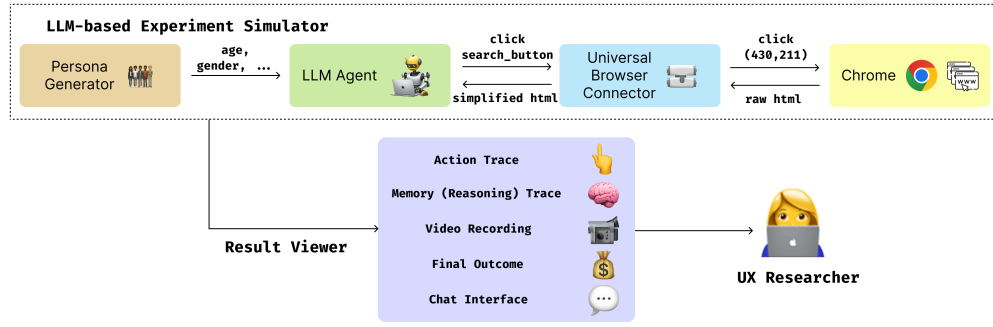
**Figure 3: System Architecture of UXAgent [45]**

Building on these advances, the use of LLM agents for the evaluation of user experience also shows promise in various domains, [63, 79, 81, 88], including graphical user interface testing [16], game environments [17, 75], and accessibility evaluation [76, 100]. For instance, SimUser [88] introduced a dual-agent framework in which one agent simulated user behaviors and another emulated the mobile app interface. While this approach showed promise for generating behavioral data automatically, it was limited by its reliance on simulated environments. As shown in Fig. 3, UXAgent [45] introduced a framework where LLM agents simulate users performing tasks on web pages, enabling automated heuristic evaluations. The persona generator generates diverse agent personas at scale, and the LLM Agent interacts with Web Browser through Web Browser Connector and produces usability data for the UX Researcher to analysis. Building on this, AgentA/B [79] applies simulated users in A/B testing contexts, demonstrating that LLM agents can detect usability improvements across interface variants in a reproducible and scalable manner.

In summary, existing approaches still lack grounding in empirical, ground-truth data derived from diverse real human participants. To address this gap, we first conducted a human-subject study to evaluate their user experiences with existing conversational shopping assistants. And conduct an agent simulation based on a real human user study.

## 3 Method

Our study consisted of two stages: (1) a ***human study*** where we invited human participants to completed two shopping tasks with conversational shopping assistants and conducted user evaluations on its usability, engagement, trust, etc and (2) an ***role play agent simulation*** that created a digital twin agent of each participant and performed the same task procedure and evaluation procedure.

In Stage 1, we recruited 40 participants to complete two shopping tasks with conversational shopping assistants. We collected demographic and background information to design realistic personas, and logged 80 shopping sessions. These sessions captured interaction traces and usability evaluations (e.g., satisfaction, task success, and perceived helpfulness), establishing a baseline for comparison with agent simulations.

In Stage 2, we scaled this evaluation through an Agent-as-a-Judge simulation. Using personas obtained from the human study, LLM-based agents role-played as digital twins for each participant,

where they completed the same tasks and provided evaluations. This allowed us to assess how closely simulated agents can reproduce human behaviors and UX judgments.

Specifically, we aim at answering three research questions (RQ) in our two-stage study:

**RQ1.** How do human customers interact with and evaluate conversational shopping assistants (CSAs)?

**RQ2.** To what extent can LLM agents role-play as customers when performing shopping tasks?

**RQ3.** How closely do agent-based customer simulations align with human behaviors in task outcomes, interaction patterns, and user experience evaluations?

### 3.1 Stage 1 Human Study: Establish Ground Truth

*3.1.1* ***Participants***. We recruited 50 participants from Prolific, aiming for a diverse U.S. sample across age, gender, race, region, education, and income level. Of these, 40 participants completed all tasks and surveys, and their data were included in our analysis. As part of a pre-screening survey, potential participants were asked to confirm that they had used an online shopping platform before, such as Amazon or Google Shopping. All studies are conducted remotely, where participants follow an instruction page that is designed to guide them through the 4-stage procedure. All participants are compensated by $15.5/h. The study protocol received ethical approval from the university's Institutional Review Board.

*3.1.2* ***Study Context - Amazon Rufus***. We conducted a user study in which participants interacted with and evaluated an agentic-AI-based conversational shopping assistant—Amazon Rufus [3] . Amazon Rufus is designed to support product discovery and decision-making through natural language dialogue. It allows users to describe their needs in free-form text, ask follow-up questions, and receive tailored product recommendations, as shown in Fig 2. Beyond simple keyword search, Rufus can summarize product features, compare alternatives, and iteratively refine suggestions across multiple conversational turns. Integrated directly into the Amazon shopping platform, it enables seamless transitions from exploration to purchase within the same web shopping interface. We selected

---

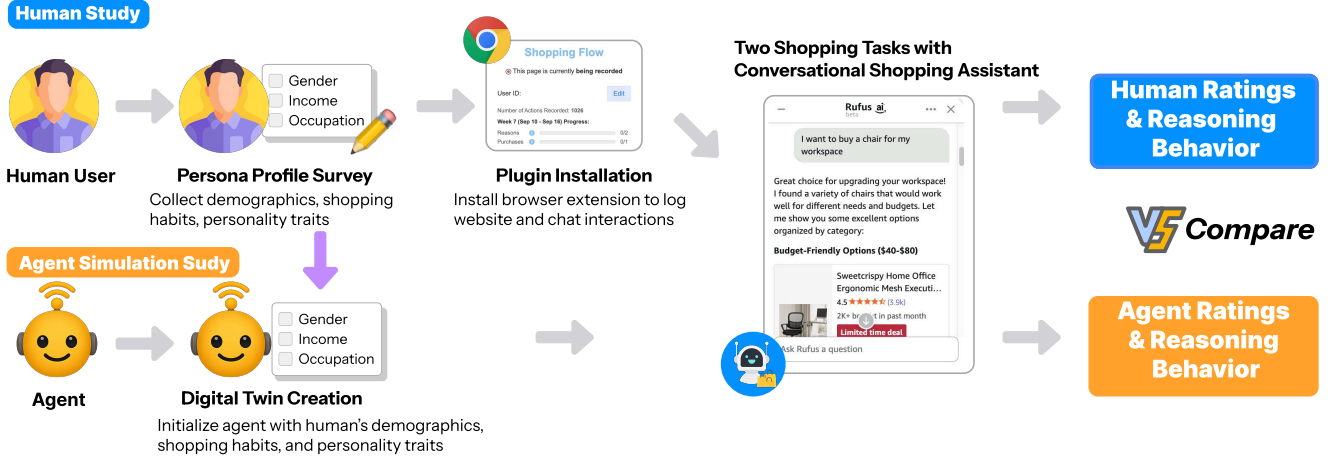[3]https://www.aboutamazon.com/news/retail/amazon-rufus

**Figure 4: Study Procedure Overview: We first collected human participants' demographics, shopping habits, and personality traits, then logged their interactions while they performed shopping tasks with a CSA. Next, we injected the same personas into LLM-based agents to complete the same tasks, and compared their ratings and reasoning behaviors to assess human–agent alignment.**

Amazon Rufus for this study because it is one of the most widely deployed conversational shopping assistants in commercial use today and offers a highly representative setting for studying user behaviors, expectations, and perceptions in shopping scenarios.

*3.1.3* ***Procedure***. The average study duration is 37.5 minutes. The overall procedure is illustrated in Figure 4.

*Step 1: Demographic Survey.* Upon entering the study, participants first completed a 15-minute structured online survey designed to elicit detailed persona profiles. The survey included items adapted from prior work shown to correlate with consumer behavior [29, 82]. It consisted of three main sections: demographic information, shopping preferences, and personality traits. Demographics covered age, gender, ethnicity and race, education, occupation, household income, and residence [29]. Shopping preferences covered topics such as online shopping frequency [71], membership status, shopping habits, seasonality [56], trust in advertising [71], engagement with product reviews, influence of delivery options [6], and an adapted eight-item Consumer Styles Inventory (CSI) [54]. Personality traits were measured using the Big Five Inventory [20] and a self-reported Myers-Briggs Type Indicator (MBTI) [53]. To contextualize each persona's familiarity with the shopping tasks, we also collected self-reported experience levels across the four product categories: monitors, chairs, outfits, and jackets. Finally, to further enrich the persona profiles, open-ended questions invited participants to describe their daily routines, shopping habits, and consumer identity[58]. Full survey items are in the Appendix.

*Step 2: Plugin Install.* Participants were asked to install a custom Chrome extension designed to capture their interactions with both the webpage and the conversational shopping assistants (CSAs). To support detailed interaction logging, we adapted the ShoppingFlow framework [82] extended the browser plugin that automatically records user behaviors and contextual signals during shopping

**Table 1: Examples of in-situ pop-up questions prompting participants to explain their interaction choices.**

| Participant Action | Pop-up Questions |
|---|---|
| *click* on related questions | Why do you want to ask rufus chatbot on user questions? |
| *click* on recommended products | We noticed that you click on product. Why did you do this? |
| *type* in the chat input | You typed in the chat input. What do you want to ask for help? |

sessions not only on the Amazon search page, but also with the Amazon Rufus Chatbot. The extension tracks granular actions such as text inputs, clicks, scrolls while interacting with Rufus, and on Amazon.com, as well as browser-level events including navigation, tab switching, and page reloads. Specifically, we focused on users' behaviors with Rufus, including typing and responding with Rufus chatbot, clicking on related questions, and clicking on recommended product, as shown in Figure 2. To better understand users' decision-making processes, the extension also featured an in-situ reflection mechanism: when participants click on related questions and recommended products, participants were prompted with a lightweight pop-up asking them to briefly explain their reasoning behind specific interactions, as shown in table 1.

*Step 3: Conduct Two Hypothetical Shopping Tasks.* Participants then completed shopping sessions using the Amazon Rufus chatbot across four tasks designed to represent both utilitarian (buying a monitor or a chair for home use) and hedonic (choosing an outfit for a summer wedding or a jacket for a group hike) shopping motivations (see Table 2). To control for task and product domain effects, participants were randomly assigned to one of two groups.

**Table 2: Participant shopping tasks of two groups. Each task specifies a category, a concrete goal, and detailed instructions to be completed. Each group completed an utilitarian task and a hedonic task.**

| Category | Task | Description |
|---|---|---|
| Utilitarian task | Find a monitor for home use | Use Rufus Chatbot to pick a $800 chair that is ergonomic for long work hours for your home workspace and add it to your shopping cart. |
| Hedonic task | Find an outfit for a summer wedding | You've been invited to a hiking event with friends. Use Rufus Chatbot to to explore hiking jackets within a $200 budget that is waterproof. If you find a desired product, add it to your shopping cart to show your intent to purchase. |
| Utilitarian task | Find a chair for home use | Use Rufus Chatbot pick a $400 monitor with the highest possible resolution for your home workspace. |
| Hedonic task | Find a jacket for a group hike | You've been invited to a summer wedding with a green theme. Use Rufus to explore your outfit options within a $200 budget. If you find a desired product, add it to your shopping cart to show your intent to purchase. |

One group completed tasks involving the purchase of a monitor and selection of a wedding outfit, while the other group completed tasks focused on purchasing a chair and selecting a hiking jacket. This grouping ensured coverage across distinct shopping goals and enabled comparison across different product categories.

*Step 4: Finish UX Evaluation Survey.* After the study session, participants filled out a survey assessing satisfaction, perceived utility, and overall experience with the assistant. After completing each task, participants filled out a user evaluation survey assessing their satisfaction, perceived usability, engagement, cognitive effort, trust, and overall experience with the shopping assistant. Appendix lists the full user evaluation survey items.

*3.1.4* **Plugin System Development**. To support detailed interaction logging, we adapted the ShoppingFlow framework [82] and developed a custom browser plugin that automatically records user behaviors and contextual signals during shopping sessions with Amazon Rufus. The plugin consists of two main components: a Content Script and a Background Script. The Content Script operates within the foreground of the active webpage, capturing fine-grained user interactions such as clicks, text inputs, scrolls, and DOM content changes. Each interaction is time-stamped and annotated with contextual metadata, including CSS selectors, semantic labels (e.g., `search_result.product_title`), and DOM attributes to enable precise mapping of actions to interface elements. The Background Script complements this by monitoring higher-level browser events

such as page navigations, tab switches, and reloads. Together, this dual-script architecture ensures comprehensive coverage of both micro-level interactions and macro-level session dynamics, while maintaining efficient performance. To better understand participants' decision-making processes, the plugin also includes an in-situ reflection mechanism. When users interact with certain UI elements—such as clicking on recommended products or related chatbot questions—the system triggers lightweight pop-up prompts asking them to briefly explain their rationale for the action.

All captured data—including user action traces, DOM snapshots, simplified HTML, and rationale responses—were stored and uploaded in real time to a secure Amazon S3 bucket whenever a user-triggered event occurred (e.g., clicking a product, typing a query, navigating to a new page). To safeguard user privacy, the plugin was explicitly configured to exclude personally identifiable information (PII) by skipping sensitive pages such as login, profile, and checkout flows. Additionally, we implemented a rule-based automated script that detects and masks any residual PII—such as usernames embedded in navigation bars, zip codes, addresses, or payment details. This instrumentation enabled the construction of a rich, multimodal dataset comprising timestamped user actions, contextual web observations, rationale annotations, and detailed persona metadata.

*3.1.5* **Measurements**. To assess both human and agent-based evaluations of conversational shopping assistants, we developed a three-part evaluation framework covering task outcome, interaction quality, and user experience. These metrics were designed to capture not only whether users completed the shopping task successfully, but also how they interacted with the assistant and how they perceived the experience. Human evaluation data were collected through post-task surveys and system-logged behavioral traces. Simulated agents were evaluated using analogous interaction logs and inferred signals. These items captured key dimensions of user experience, including task success, user satisfaction, and user perceptions of the interaction—such as product information accuracy, usability, efficiency, conversational quality, trust, and cognitive load.

*Shopping Task Outcomes.* We measured task success based on a combination of user self-reports and final product selections. After each shopping session, participants confirmed whether they used the assigned assistant (i.e., Amazon Rufus) and indicated whether they were able to find a product that matched the given shopping goal. Participants also rated their satisfaction with the final outcome on a 5-point Likert scale, providing a subjective measure of task success.

*User Interaction Data.* Interaction quality was assessed using a combination of automatically logged behavioral data from the Chrome extension and self-reported user feedback. From the interaction logs, we analyzed the number of clarification queries made during each session, reflecting how actively users engaged in refining their requests and how well the assistant supported iterative dialogue. We also examined the depth and specificity of user queries, particularly when users asked about nuanced product features, comparisons, or trade-offs—indicating the extent to which the assistant enabled informed decision-making.

*User Experience Survey.* Participants completed a short post-task survey rating their experience with the CSAs. We captured their evaluations on usability, engagement, satisfaction, trust, and cognitive effort. Usability was captured using adapted items from the System Usability Scale (SUS) [8], such as "It was easy to interact with Rufus" and "I found the interaction enjoyable." Satisfaction and intention to reuse were measured using items like "I would love to use Rufus to shop in the future" and "I will recommend others to use Rufus to shop." Perceived helpfulness was assessed through statements such as "Rufus helped me narrow down product options." Conversational engagement was also measured with the statement "The conversation with Rufus felt engaging."

To complement behavioral data, participants also self-reported their engagement levels and cognitive effort [13, 25]. Example statements included "I found the interaction mentally demanding" and "I had to put in a lot of effort to use Rufus effectively," both adapted from the NASA Task Load Index (NASA-TLX)[25] and the Technology Acceptance Model (TAM)[13].

In addition, participants evaluated the assistant's information accuracy and performance [5]. These included items such as "Rufus provided accurate and up-to-date product information," "Rufus's responses were logical and coherent," and "Rufus described product features in a way that matched official information," which collectively measured perceived information accuracy, coherence, and factual consistency.

We also included measures of perceived trustworthiness and reliance, which are essential in assessing user confidence in AI recommendations. Participants rated statements such as "I trust the responses provided by Rufus," "I would rely on Rufus without double-checking its responses," and a reverse-coded item: "I was concerned that Rufus may present biased or sponsored recommendations." These items reflect emerging concerns around AI bias, transparency, and user trust in conversational agents [5].

At the end of the survey, participants were asked an open-ended question: "What do you like about Rufus? What do you dislike?" This allowed us to capture qualitative insights into user-perceived strengths, weaknesses, and unmet expectations that may not be captured by structured response scales. Together, these measures provide a multidimensional view of user experience across usability, utility, trust, cognitive effort, and emotional response, enabling both quantitative comparison and in-depth qualitative interpretation.

*3.1.6* **Data Analysis**. To establish a human benchmark, we analyzed data from 40 participants who completed four shopping tasks (monitor, chair, summer outfit, hiking jacket), spanning both utilitarian and hedonic goals. Task success was binary-coded based on whether the selected product met task-specific constraints (e.g., staying within budget, suitability for the occasion). Post-task user experience (UX) ratings were collected on five dimensions—satisfaction, trust, usability, helpfulness, and cognitive load—using 5-point Likert scales (1 = strongly disagree, 5 = strongly agree).

We computed descriptive statistics and conducted paired and independent-sample $t$-tests to compare task outcomes and UX ratings across task types (utilitarian vs. hedonic). Interaction behavior metrics—including total turn count, number of clarification queries, and average message length—were extracted from logged sessions and analyzed.

## 3.2 Stage 2 LLM Agent Simulation: Role-Playing as Digital Twins

In Stage 2, we extended the formative human study by scaling the evaluation using an Agent-as-a-Judge simulation approach. Building on real participant data, we instantiated LLM-based agents conditioned on participant personas, prompting them to role-play as diverse online shoppers. These digital twin agents completed the same shopping tasks as human participants and generated UX evaluations, enabling large-scale, repeatable assessments of conversational shopping assistants—grounded in real-world behavioral patterns.

To simulate human experiences and evaluate AI-powered shopping interactions at scale, we employed the UXAgent framework [45]. UXAgent is a persona-driven LLM agent architecture designed to model user behavior through realistic task execution and reflective UX reasoning. In our study, UXAgent was used to simulate **"digital twins"** of real participants from the human study, enabling direct comparisons between human and agent performance and perception under identical task conditions. We selected UXAgent because of its explicit design for planning and reflection during UX evaluations, which made it well suited to our study goals.

*3.2.1* **Transform from Real Participants to Persona**. We grounded our agent simulations in the behavioral, contextual, and evaluative patterns observed in the human study (Section 3.1). Using demographic profiles, shopping goals, interaction traces, and post-task feedback from 50 participants, we constructed a one-to-one mapping between each human user and their simulated digital twin. Each agent was instantiated with a benchmark persona that combined demographic and contextual attributes—such as age group, gender, shopping frequency, tech familiarity, and budget sensitivity—with shopping styles and priorities inferred from observed behavior, including emphasis on price, brand loyalty, reliance on product comparisons, and the use of clarification strategies. To enrich the persona beyond static attributes, we also incorporated participants' self-descriptions and daily routines, collected through open-ended survey questions. The resulting persona descriptions served as natural language prompts to guide agent behavior, ensuring that each "digital twin" reflected not only the participant's demographic background and shopping logic, but also their lived experience and decision-making tendency.

*3.2.2* **Procedure**. The simulation followed four core stages:

*Persona Initialization.* Each agent was primed with a natural language persona prompt summarizing the participant's background, shopping motivations, and task-specific constraints. For example: *"You are a 34-year-old professional who prefers eco-friendly products under $400. You are shopping for a monitor for your home office. You care about display quality and customer reviews."* This setup enabled the agent to behave consistently with the user's profile throughout the task.

*Plan Generation and Task Execution.* The agent then begins task execution in an iterative loop of planning and acting. In each iteration, it generates or revises a structured, step-by-step plan to complete the assigned shopping task (e.g., product search, comparison,

clarification, selection), grounded in the persona's goals and constraints. The agent then executes this plan via UXAgent's *Universal Web Connector* [45], which enables operation in real-world web environments (including Amazon Rufus) using simplified HTML observations and task-agnostic action primitives.

The Web Connector filters raw HTML into a semantic representation that preserves meaningful interface elements such as product titles, chatbot responses, buttons, and related questions, while discarding visual clutter (CSS, JavaScript). Based on this abstraction, the agent selects from a fixed set of atomic actions: `click`, `type`, `type_and_submit`, `clear`, `back`, and `terminate`. These actions are executed through a backend to simulate realistic web interactions.

Each action prediction was made using the Claude 3.7 Sonnet model, selected for its strong reasoning and instruction-following capabilities. The model operated under a low-temperature decoding setting (`temperature=0.2`) to prioritize consistency and reproducibility of behavior. For each step, the model received the persona context, current task goal, interaction history, and simplified DOM snapshot as input, and predicted the next best action. This dynamic reasoning loop—plan, act, observe, revise—enabled agents to adapt to unexpected responses, revise queries, or switch strategies mid-task, mimicking realistic user behavior.

*Post UX Evaluation.* Upon completing each shopping task, the agent was prompted to evaluate the assistant's performance using their interaction tracing. We asked agents to answer the same survey questions as users. The agent answers the open-ended questions and rates the Likert-style ratings (1–5) across the same dimensions used in the human study: task success, usability, conversational quality, trust, and cognitive load. These role-aligned self-assessments leveraged UXAgent's prompting strategy to ensure consistent and interpretable judgments.

*Interaction and Session Logging.* All interaction data, including the agent's queries, Rufus's responses, selected actions, clicked elements, and internal planning steps, were logged in detail. We also stored the action traces and final product selections. This enabled fine-grained comparison with human users.

This end-to-end simulation pipeline enabled controlled, persona-consistent evaluation of conversational shopping assistants at scale. By integrating UXAgent's reasoning architecture with its Universal Web Connector for real-world web execution, we were able to simulate human-like shopping behavior on Amazon Rufus and evaluate agentic alignment with real user goals, preferences, and perceptions.

*3.2.3* ***Measures and Analysis***. To evaluate how well persona-grounded LLM agents align with human users, we conducted a multi-level comparative analysis grounded in the same dimensions assessed in our human study. Simulated agent behaviors were evaluated using both structured outputs (e.g., UX ratings generated at the end of each task) and detailed interaction logs that are automatically captured during the simulation. The analysis was designed to support both qualitative comparison and statistical testing of alignment between humans and their digital twin agents.

*Persona-Level Alignment.* We assessed whether the agent's final product choices were consistent with the preferences encoded in its assigned persona. Each agent was grounded in a real participant's demographic profile, shopping goals, and stated priorities. We compared final product selections between the agent and the corresponding human participant to measure choice consistency.

*Task Outcome Alignment.* To evaluate task success, we calculated the success rate of purchasing products. Then, we computed the F1 score on their decision on purchase or not.

*Interaction Behavior Alignment.* We examined behavioral measures extracted from simulation logs, including the number of dialogue turns, number of recommendation products, and number of related questions. These metrics were directly comparable to human interaction traces. We used Welch's t-test to compare human with their paired agent. We also compared the semantic of messages including the length of messages and cosine similarity of messages. To compare the action trajectory between agent and human, we further calculated the Levenshtein distance of their action sequences. This analysis allowed us to capture both structural alignment and strategic divergence.

*UX Evaluation Alignment.* At the end of each session, agents generated structured UX evaluations using prompts aligned with the human post-task survey. Ratings were captured across the same five dimensions: task satisfaction, perceived helpfulness, usability, conversational quality, and cognitive load. While human ratings were self-reported, agent ratings were inferred via reflective prompting based on the agent's persona and full action trace. To compare these ratings, we used Welch's t-test on the Likert-scale responses, quantifying overestimation tendencies (e.g., agents consistently rating satisfaction higher).

## 4 Result

We analyzed 40 participants' data across 80 human shopping sessions and their digital twins' data across 80 agent-simulated sessions to answer research questions.

### 4.1 RQ1: How do human customers interact with and evaluate CSAs?

*4.1.1 Task Outcomes.* All participants successfully purchased related items, with 23% selecting the same products. Reported satisfaction with chosen items was high (M = 4.5, SD = 0.91 on a 5-point scale). The average shopping time was 375.9 seconds (SD = 203). A Kolmogorov–Smirnov test indicated no significant distributional differences between the group that selected chair and jacket vs. the gorup that selected monitor and outfit, in terms of number of dialogue turns and shopping time. Therefore, we combined the two groups in subsequent analyses (N = 40). Compared with two types of shopping tasks, which is a utilitarian task versus a hedonic task, participants spent more turns on hedonic tasks (M=2.1, SD=1.3) than utilitarian tasks (M=1.8, SD=1.2) to explore potential options.

*4.1.2 Interaction Traces.* Analysis of human interaction behaviors provides insight into how participants engaged with the conversational shopping assistants. On average, participants spent 341.1 seconds (SD = 182.7) completing a shopping session. Sessions contained an average of 27.3 interaction events (SD = 11.8), including 11.7 clicks (SD = 5.3) and 4.8 typed inputs (SD = 2.9), suggesting

that users balanced passive browsing of recommendations with active search queries.

Message-level analysis further revealed that participants contributed an average of 1.9 customer messages per session (SD = 1.2), with most users sending between one and two utterances. These messages typically focused on clarifying constraints (e.g., budget, product features) or probing for alternatives before making a final choice. One example of a user chat message with Rufus is shown in Fig. 5

*4.1.3 UX Evaluations.* Participants' narratives revealed a tension between appreciating Rufus's efficiency and noticing its limitations. Many praised its ability to quickly surface relevant products, describing it as "helpful in showing me a list of many things that I was searching for [P5]" and "supplying options that were exactly what I asked for [P7]." Some even compared the experience to having a personal assistant, with one participant noting that "it's like having a personal shopper! [P11]" These accounts underscore Rufus's strength in streamlining product discovery, especially when users had clear specifications in mind.

At the same time, users also highlighted moments of frustration and comparison with manual browsing. Several mentioned being shown inappropriate or mismatched items, such as one participant who felt "frustrated since it got my gender wrong and was trying to show me male clothing at first." Others pointed out the loss of breadth compared to traditional search: "I can't see as many options at a time when using Rufus compared to manually searching" and "manual search allows me to see more options and compare products." Looking ahead, some participants expressed enthusiasm for integrating AI assistants into their shopping routine—"I will definitely use Rufus in the future to help me begin my searches"—while others emphasized conditional adoption, saying they would only use it if recommendations became more transparent and context-aware. Together, these reflections highlight both the promise of Agentic-based-AI shopping assistants in reducing effort and the need to address personalization gaps and user control.

## 4.2 RQ2: Can LLM agents role-play as customers when performing shopping tasks?

*4.2.1 Task Outcomes.* Agents were able to complete the assigned shopping tasks, but their product choices diverged from those made by humans. Direct overlap between the two groups was rare—only 1.3% of cases involved both an agent and a human participant selecting the exact same product. For the hedonic tasks that the agent can make a decision on whether to purchase or not. Agents decided to add product to the shopping carts in 76 out of 80 sessions.

*4.2.2 Interaction traces.* Analysis of agent interaction traces shows that simulated customers engaged in concise but systematic exchanges with the shopping assistant. On average, agents produced 2.1 customer messages per session (SD = 1.3), and these messages tended to be longer in character length (M = 218, SD = 123 across all messages). Agents consistently started with well-structured first queries averaging 7.8 words (SD = 2.3), often explicitly stating product constraints such as category, budget, or features. Across sessions, agents executed an average of 10.3 total actions (SD = 5.2), including 5.6 clicks (SD = 3.7) and 3.5 typed inputs (SD = 1.8),

illustrating a balanced strategy of browsing recommendations and issuing directed prompts.

Beyond overall action counts, click behavior highlights the agents' tendency toward broader exploration. On average, agents clicked on 1.9 recommended items (SD = 1.3), while also occasionally engaging with related questions (M = 0.78, SD = 0.8). This frequency of interactions with multiple recommendation types suggests that agents actively probe the system for alternatives before finalizing a decision. Despite this exploratory behavior, every agent session included a successful add-to-cart action, demonstrating that agents were able to complete assigned shopping tasks reliably. In post-task reflections, agents reported a mean self-satisfaction rating of 3.97 (SD = 0.59) on a five-point scale, indicating consistently positive perceptions of their choices. Taken together, these traces reveal that simulated agents mirror humans in structural metrics but adopt a more exhaustive and systematic exploration pattern, while also displaying a tendency to evaluate outcomes more positively than human participants.

*4.2.3 UX Evaluations.* The majority of agent-generated responses emphasized Rufus's helpfulness, efficiency, and organization. Across tasks, agents frequently described Rufus as providing "curated, categorized recommendations" or "organizing monitors by resolution categories," highlighting a consistent framing of the assistant as a structured and efficient tool. These responses suggest that simulated agents tend to evaluate interactions in narrowly functional terms, focusing on speed and clarity of information delivery. Unlike human participants, they rarely reported frustration or ambiguity, instead portraying Rufus as reliably helpful. This pattern illustrates how agent role-play produces a more uniformly positive account of the shopping experience.

At the same time, agents often framed their feedback in direct comparison to traditional search or browsing. Many responses claimed that Rufus "organized information more efficiently than manual search[A8]" or "provided better product matches than traditional browsing[A4]" While human participants sometimes highlighted frustrations with reduced control or missing context, the agents did not mention such limitations. Finally, a smaller set of responses reflected future intentions, with agents stating they would "use AI shopping assistants more frequently [A12]" or "increase reliance on Rufus for product discovery [A33]" Taken together, these findings show that agent responses gravitate toward highlighting efficiency and structure, with less diversity and nuance than human narratives.

## 4.3 RQ3: How closely do agents align with humans in task outcomes, interaction patterns, and user experience evaluations?

*4.3.1 Task Outcomes: Agents Complete Tasks but Choose Differently from Humans.* Agents were able to complete the assigned shopping tasks, but their product choices diverged from those made by humans. Agent matches with humans on their final buy-or-not decisions (F1 score of 0.9). Among agents, 45% of selections were duplicates, compared to 23% among humans, indicating that agents converged on the same items slightly more often. However, direct overlap between the two groups was rare—only 1.3% of cases
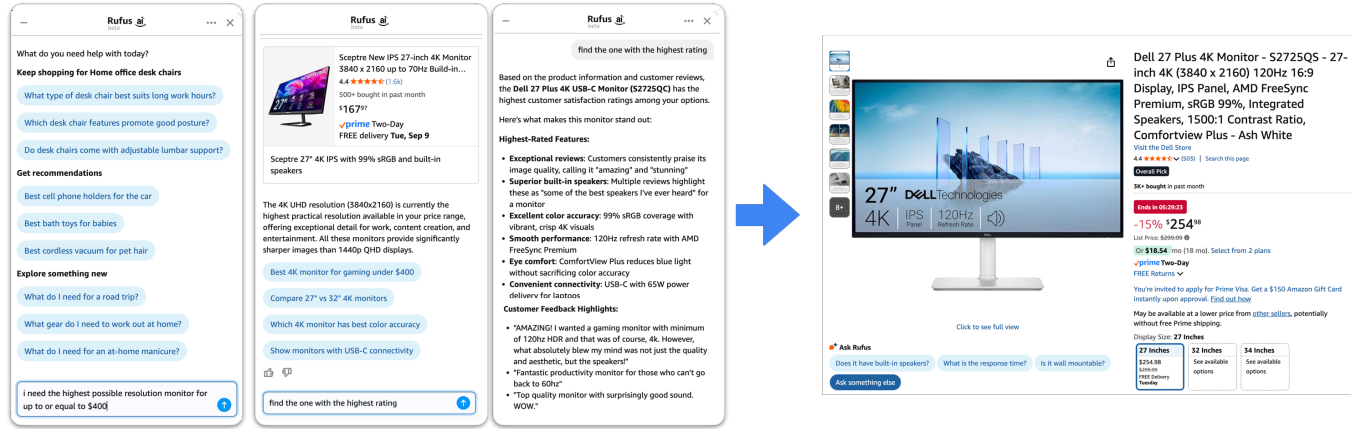
**Figure 5: Example of human participant chatting with Rufus to buy a monitor that is under** $400 **with the highest resolution based on their own preferences**

involved both an agent and a human participant selecting the exact same product. For example, participant [P18] selected the "LG 27UP850K-W 27-inch Ultrafine 4K UHD", while the corresponding agent [A18] selected the product "MSI PRO MP273U, IPS 3840 x 2160 (UHD) Computer Monitor, 4K". Furthermore, participants reported being more satisfied with their product selections than the agents, suggesting that while agents successfully completed tasks, their choices did not align as closely with human preferences. Based on the results, agents can complete tasks effectively. However, their decision trajectories differ from those of humans—likely due to the agent's use of step-by-step action prediction, where each decision is tightly coupled to the previous one [84], in contrast to humans' heuristics, preferences, or cognitive biases in their final selections [77].

*4.3.2 Interaction traces: Agents and humans start alike but explore differently.* Analysis of interaction traces revealed both structural similarities and behavioral divergences between humans and agents. The number of dialogue turns per session did not differ significantly between groups (humans: $M = 1.9$, $SD = 1.2$; agents: $M = 2.1$, $SD = 1.3$), suggesting that agents can replicate the high-level pacing of human interactions. Both humans and agents also tended to initiate conversations with similarly structured first messages, typically stating the product type and budget constraint—indicating successful alignment in initial task framing. The cosine similarity of the first message using sentence embeddings is average 0.49, which indicates that they shared some similar wording in the first round messages. We transformed human and agent action traces into trajectory sequences (e.g., [ask_rufus_a_question, click_related_question, ask_follow_up_question, click_product, click_review, add_to_cart]) and compared them using Levenshtein distance. The normalized edit distance (NED = 0.89, SD = 0.06) suggests substantial divergence between human and agent trajectories, while the similarity score (SIM = 0.11, SD = 0.06) indicates only limited overlap, underscoring that agents often follow systematically different action paths than humans. This divergence reflects how agents tend to explore more broadly

and exhaustively, whereas humans rely on selective, goal-directed strategies to reach decisions more efficiently.
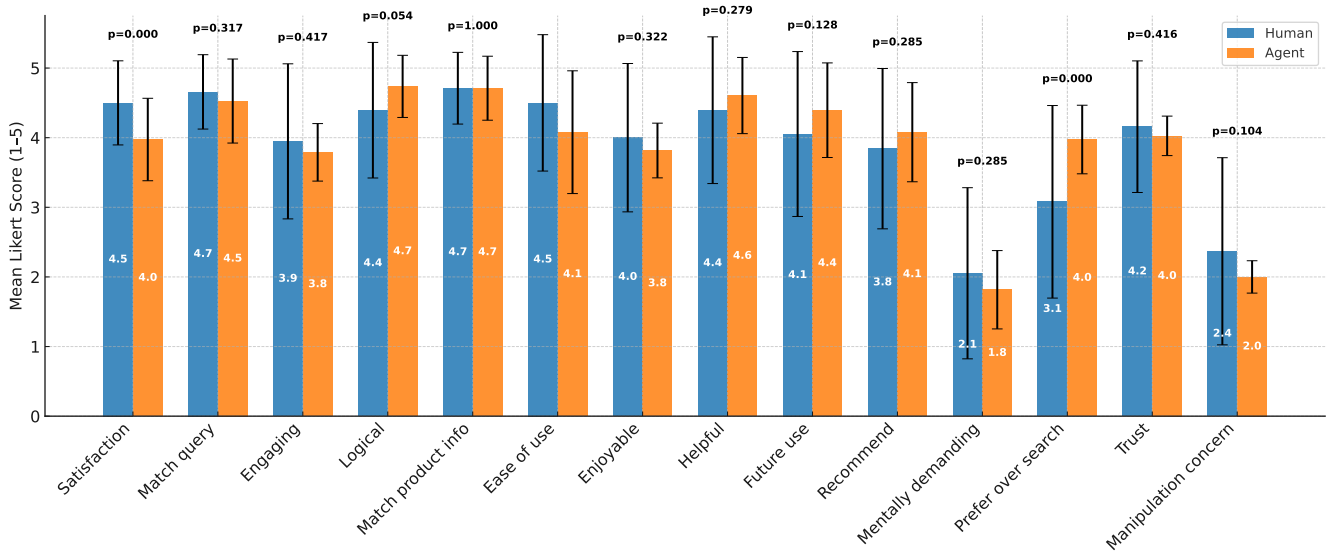
This contrast is reflected in several quantitative measures. Agents clicked on significantly more recommended items ($M = 1.9$, $SD = 1.3$) than humans ($M = 1.2$, $SD = 1.0$, Welch's t-test $p < 0.01$). They also asked more related follow-up questions ($M = 0.78$, $SD = 0.8$) compared to humans ($M = 0.13$, $SD = 0.7$). In terms of message semantics, agents' first messages were significantly shorter in tokens ($M = 7.8$, $SD = 2.3$) than those of humans ($M = 14.8$, $SD = 4.7$, Welch's t-test $p < 0.001$), yet their cumulative message length was much greater ($M = 218$, $SD = 123$ vs. $M = 120$, $SD = 60$) in character, reflecting a more exploratory dialogue style.

Interestingly, not all interaction measures differed significantly. For example, the number of dialogue turns did not show significant group differences and the cosine smiliarity of the first messages is higher than 0.4, suggesting that both groups engaged in similarly structured conversation lengths. Together, these findings indicate that while agents mirror the rhythm and opening structure of human conversations, their internal decision processes emphasize breadth-first exploration, surfacing a wider range of product options than humans typically consider.

*4.3.3 UX evaluations: Agents tend to prefer Rufus over the traditional interface, while humans are more satisfied with the product output.* We compared post-task UX survey ratings between human participants and their simulated customer agents using paired t-tests. Overall, ratings were broadly consistent across most dimensions, with no significant differences in perceived query matching, coherence, enjoyment, helpfulness, or trust (all $p > 0.1$). Two dimensions showed significant differences: humans reported higher satisfaction with the final product they picked than their simulated digital twins (Welch's t-test $p < 0.001 * **$), and they were more likely to prefer Rufus over traditional search (Welch's t-test $p < 0.001 * **$). These results suggest that while simulated agents can approximate human evaluations on most UX aspects, they systematically underrate satisfaction with the chosen outcome and the comparative advantage of conversational shopping over traditional search.

**Table 3: Comparison of interaction measures between humans and agents. Asterisks indicate statistical significance on t-test (* $p < 0.05$, ** $p < 0.01$),*** $p < 0.001$).**

| Interactions | Measure | Humans | Agents | |
|---|---|---|---|---|
| Actions | # of turns | 1.9 (SD=1.2) | 2.1 (SD=1.3) | |
| | # of recommended items | 1.2 (SD=1.0) | 1.9 (SD=1.3) | ** |
| | # of related questions | 0.13(SD=0.7) | 0.78(SD=0.8) | ** |
| Message semantic | Length of first message | 14.8 (SD=4.7) | 7.8(SD=2.3) | *** |
| | Length of all messages | 120 (SD=60) | 218 (SD=123) | *** |
| | Cosine similarity of the first message | 0.49 (SD=0.32) | | |
| | Normalized Levenshtein distance | NED: 0.89 (SD=0.06), SIM: 0.11(SD=0.06) | | |



**Figure 6: Post-Study UX Survey Result Comparison.**

Both humans and simulated customer agents consistently highlighted Rufus's strengths in efficiency and product organization. Participants described Rufus as "helpful in showing me a list of many things that I was searching for [P23]," while agents echoed similar language, calling it "curated" and "organized more efficiently than manual search [P2]". This overlap suggests that both groups recognize and value Rufus's ability to streamline the shopping process, particularly for utilitarian tasks where speed and relevance matter. In this respect, agent role-play captures some of the functional aspects of human evaluation.

However, notable differences emerged in how each group articulated limitations and subjective experiences. Human participants frequently reported frustrations or gaps in personalization—such as "I felt frustrated when Rufus thought I was a woman when I was looking for the green themed wedding outfit" or "I can't see as many options at a time when using Rufus compared to manually searching.[P31]" In contrast, agent-generated responses rarely mentioned negative experiences, instead focusing almost exclusively on structured comparisons to search engines (e.g., "Rufus organized information more efficiently than traditional search [A4]").

Humans also expressed ambivalence about future use, with some welcoming AI shopping support and others warning against manipulation, whereas agents framed future intentions more uniformly as increased reliance. Together, these findings show that while agents reproduce positive evaluations of efficiency, they lack the nuance and critical perspective present in human accounts, especially around personalization, trust, and control.

## 5 Discussion

In this study, we used agents to judge agentic-AI systems in the context of conversational shopping assistants. Our study is the first to quantify how well LLM agents can simulate specific human users in multi-turn, goal-driven shopping tasks. While prior work has examined LLMs as judges or single-turn role-players, our results provide the first empirical evidence of alignment with human benchmarks in real-world shopping interactions. The findings highlight both the promise and the limits of current LLM-based simulations. On the one hand, agents aligned closely with humans on structural measures such as task completion, purchase decisions, dialogue length, and coherent evaluations. Like human participants, they

consistently recognized Rufus's efficiency and organization, suggesting that agents can serve as scalable proxies for benchmarking system performance, particularly in early-stage evaluations where recruiting diverse participants is costly or infeasible.

At the same time, important divergences emerged. Although humans and agents opened conversations in similar ways, agents explored a broader range of product options before deciding, reflecting a systematically different decision-making style. These discrepancies underscore that agents can approximate functional judgments, but they fall short in capturing the affective and critical dimensions of human experience. Together, these findings highlight the need for hybrid evaluation strategies that combine the scalability of agents with the nuance of human judgment, as well as future work on bridging the gap in reasoning and exploration between LLM agents and real users.

## 5.1 Simulating Humans with Role-Play Agents in Multi-Turn Tasks

Role-play agents reproduced several structural aspects of human interaction, including similar turn counts, coherent first queries, and consistent task completions. Yet their trajectories quickly diverged. Agents favored breadth-first exploration—clicking more recommendations and issuing longer queries—whereas humans narrowed rapidly to decision-relevant constraints. This suggests that while current LLM agents can approximate the form of interaction, they do not yet capture the underlying reasoning processes that guide human decision-making in adaptive, longitudinal tasks [95].

## 5.2 Bridging the Gap in Reasoning and Exploration

The divergence between human and agent trajectories highlights an important research opportunity. Humans rely on heuristics, prior experience, and bounded rationality, while agents pursue more systematic but less preference-sensitive strategies. This mismatch produced only 2% overlap in exact product choices, despite comparable task completion rates. Future work could focus on training or fine-tuning LLM agents with human-like heuristics, decision biases, and preference models to close the gap in reasoning and exploration between digital twins and real users [43].

## 5.3 Agent-as-a-Judge for UX Simulation Studies

Our findings suggest several implications for the design and evaluation of CSAs. First, LLM agents can serve as scalable proxies for early-stage evaluation, providing quick feedback on task success, coherence, and recommendation quality without the cost of recruiting large participant samples. Prior work in the search context similarly proposed that agentic search can augment human cognition by automating routine legwork and allowing people to focus their limited capacity on higher-order concerns such as critical decisions and oversight [22]. Second, our results highlight the importance of hybrid evaluation strategies: while agents provide breadth and consistency, they fall short in capturing affective dimensions such as satisfaction and preference, making human evaluation indispensable. Third, persona-driven simulation enables researchers and practitioners to probe diverse customer scenarios that are often difficult to recruit for in small-scale studies, thereby broadening

coverage of customer experience. Finally, designers should treat agent-generated evaluations as directional signals rather than definitive judgments, particularly when making design decisions that may influence user engagement and trust.

## 5.4 Human-Agent Collaboration for Hybrid UX Evaluation

Our findings underscore the need for hybrid evaluation. LLM agents offer scalability and speed, making them well suited for early-stage assessments. However, their actions are divergent compared with human participants. Embedding humans in the loop—whether by calibrating substitution costs for agent decision-making, or by periodically anchoring simulations with human feedback will ensure that agent-based evaluations remain grounded in lived experience. Potential future work can further develop visualization or replay mechanism that help designers and researchers to further understand the agents' reasonings and thinking behavior [45].

Our goal in introducing agents into UX evaluation is not to replace human participants but to complement them. Agents can generate early signals and pilot testing, giving UX designers more opportunities to iterate quickly. Yet the divergences we observed in trajectories, satisfaction ratings, and nuanced judgments, make clear that agents cannot capture the full spectrum of human experience. Relying solely on simulations risks overlooking critical user concerns, such as issues of fairness, personalization, or trust. Keeping humans in the evaluation loop is therefore both a methodological and an ethical imperative, ensuring that design decisions are accountable to real users and that agent-based insights are treated as directional rather than definitive.

## 5.5 Simulating Human Behaviors and Societal Dynamics

Our work contributes to the broader vision of using LLM agents to model diverse populations of users. Prior research has demonstrated the potential of simulating hundreds or even thousands of agents to examine emergent behaviors in domains such as negotiation, economics, and social interaction [58, 101]. Our study extends this line of work by grounding simulations in empirical human data, enabling more faithful digital twins. Scaling from dozens of such twins to thousands of synthetic users opens opportunities to investigate systemic questions, such as trust in AI shopping assistants or marketing strategies. Looking forward, this line of research can inform not only HCI design but also policy discussions on the societal impact of agentic AI.

## 6 Limitation

While our study provides new insights into the use of LLM agents as digital twins for evaluating agentic AIs, several limitations should be noted.

First, our evaluation focused on a single domain—online shopping—and a single platform, Amazon Rufus. Although Rufus is a widely deployed system, findings may not generalize to other shopping assistants with different interaction designs, product domains, or recommendation logics. Moreover, other forms of agentic AIs, such as conversational assistants in productivity, education,

or healthcare—remain to be explored. These domains may present distinct challenges for simulating user behavior and evaluating UX.

Second, our study examined a limited set of tasks—two utilitarian (monitor, chair) and two hedonic (outfit, jacket). While these capture common shopping goals, they do not reflect the full breadth of real-world online shopping behaviors. Broader task sets will be important to examine how well agents model diverse goals, from everyday purchases to more open-ended exploration.

Finally, our agent simulations relied on one agent implementation – UXAgent – built an off-the-shelf Claude 3.7 Sonnet model. We selected UXAgent because of its explicit design for planning and reflection during UX evaluations, which made it well suited to our study goals. Nonetheless, model capabilities and alignment strongly influence agent decision-making and evaluation patterns, and different LLMs or agent architectures may yield varying levels of human–agent alignment. As LLMs and agent frameworks continue to evolve, replicating our study with alternative models and agent designs will be essential for assessing the robustness and generalizability of our findings.

## 7 Conclusion

Our study provides the first quantitative evidence of how closely LLM agents can simulate human customers in multi-turn conversational shopping tasks. We found that while agents matched humans in structural measures such as task completion and turn count, they diverged in product choices, interaction strategies, and subjective evaluations, with humans reporting higher satisfaction and trust. LLM-agent judges aligned well with human raters on objective dimensions like task success and relevance but systematically overestimated satisfaction and helpfulness. These findings highlight the value of LLM agents for scalable, early-stage evaluation of conversational shopping assistants, while underscoring the continued importance of human studies for capturing nuanced, affective aspects of user experience.

## References

[1] Deepak Bhaskar Acharya, Karthigeyan Kuppan, and B Divya. 2025. Agentic ai: Autonomous intelligence for complex goals–a comprehensive survey. *IEEe Access* (2025).

[2] Elske Ammenwerth, Stefan Gräber, Gabriele Herrmann, Thomas Bürkle, and Jochem König. 2003. Evaluation of health information systems—problems and challenges. *International journal of medical informatics* 71, 2-3 (2003), 125–135.

[3] Terence S Andre, H Rex Hartson, Steven M Belz, and Faith A McCreary. 2001. The user action framework: a reliable foundation for usability engineering support tools. *International Journal of Human-Computer Studies* 54, 1 (2001), 107–136.

[4] Daniel Baier, Alexandra Rese, Maximilian Röglinger, D Baier, A Rese, and M Röglinger. 2018. Conversational User Interfaces for Online Shops? A Categorization of Use Cases.. In *ICIS*.

[5] Janarthanan Balakrishnan and Yogesh K Dwivedi. 2024. Conversational commerce: entering the next stage of AI-powered digital assistants. *Annals of Operations Research* 333, 2 (2024), 653–687.

[6] Živilė Baubonienė and Gintarė Gulevičiūtė. 2015. E-commerce factors influencing consumers 'online shopping decision. *Social technologies* 5, 1 (2015), 62–73.

[7] Petter Bae Brandtzaeg and Asbjørn Følstad. 2018. Chatbots: changing user needs and motivations. *interactions* 25, 5 (2018), 38–43.

[8] John Brooke et al. 1996. SUS-A quick and dirty usability scale. *Usability evaluation in industry* 189, 194 (1996), 4–7.

[9] Chaoran Chen, Zhiping Zhang, Ibrahim Khalilov, Bingcan Guo, Simret A Gebreegziabher, Yanfang Ye, Ziang Xiao, Yaxing Yao, Tianshi Li, and Toby Jia-Jun Li. 2025. Toward a human-centered evaluation framework for trustworthy llm-powered gui agents. *arXiv preprint arXiv:2504.17934* (2025).

[10] Jiaju Chen, Yuxuan Lu, Xiaojie Wang, Huimin Zeng, Jing Huang, Jiri Gesi, Ying Xu, Bingsheng Yao, and Dakuo Wang. 2025. Multi-Agent-as-Judge: Aligning LLM-Agent-Based Automated Evaluation with Multi-Dimensional Human Evaluation. arXiv:2507.21028 [cs.CL] https://arxiv.org/abs/2507.21028

[11] Leigh Clark, Nadia Pantidi, Orla Cooney, Philip Doyle, Diego Garaialde, Justin Edwards, Brendan Spillane, Emer Gilmartin, Christine Murad, Cosmin Munteanu, et al. 2019. What makes a good conversation? Challenges in designing truly conversational agents. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–12.

[12] Nils Dahlbäck, Arne Jönsson, and Lars Ahrenberg. 1993. Wizard of Oz studies: why and how. In *Proceedings of the 1st international conference on Intelligent user interfaces*. 193–200.

[13] Fred D Davis. 1989. Perceived usefulness, perceived ease of use, and user acceptance of information technology. *MIS quarterly* (1989), 319–340.

[14] Steven Dow, Blair MacIntyre, Jaemin Lee, Christopher Oezbek, Jay David Bolter, and Maribeth Gandy. 2005. Wizard of Oz support throughout an iterative design process. *IEEE Pervasive Computing* 4, 4 (2005), 18–26.

[15] Shiyu Duan. 2025. Systematic analysis of user perception for interface design enhancement. *Journal of Computer Science and Software Applications* 5, 2 (2025).

[16] Juha Eskonen, Julen Kahles, and Joel Reijonen. 2020. Automating GUI Testing with Image-Based Deep Reinforcement Learning. In *2020 IEEE International Conference on Autonomic Computing and Self-Organizing Systems (ACSOS)*. 160–167. doi:10.1109/ACSOS49614.2020.00038

[17] Pedro M. Fernandes, Manuel Lopes, and Rui Prada. 2021. Agents for Automated User Experience Testing. In *2021 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*. 247–253. doi:10.1109/ICSTW52544.2021.00049

[18] Firework. 2023. What Is Conversational Shopping? https://firework.com/blog/what-is-conversational-shopping Accessed: 2025-06-05.

[19] GitHub and OpenAI. 2021. GitHub Copilot: Your AI Pair Programmer. https://github.com/features/copilot/. Accessed: 2025-09-10.

[20] Lewis R Goldberg. 1992. The development of markers for the Big-Five factor structure. *Psychological assessment* 4, 1 (1992), 26.

[21] Google DeepMind. 2023. Gemini: A Family of Highly Capable Multimodal Models. arXiv preprint. arXiv:2312.11805, Accessed: 2025-09-10.

[22] Boyu Gou, Zanming Huang, Yuting Ning, Yu Gu, Michael Lin, Weijian Qi, Andrei Kopanev, Botao Yu, Bernal Jiménez Gutiérrez, Yiheng Shu, et al. 2025. Mind2Web 2: Evaluating Agentic Search with Agent-as-a-Judge. *arXiv preprint arXiv:2506.21506* (2025).

[23] Jan Gulliksen, Bengt Göransson, Inger Boivie, Stefan Blomkvist, Jenny Persson, and Åsa Cajander. 2003. Key principles for user-centred systems design. *Behaviour and Information Technology* 22, 6 (2003), 397–409.

[24] Onder Gurcan. 2024. LLM-Augmented Agent-Based Modelling for Social Simulations: Challenges and Opportunities. doi:10.48550/arXiv.2405.06700 arXiv:2405.06700 [physics]

[25] Sandra G Hart and Lowell E Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of empirical and theoretical research. In *Advances in psychology*. Vol. 52. Elsevier, 139–183.

[26] Rex Hartson and Pardha S Pyla. 2012. *The UX Book: Process and guidelines for ensuring a quality user experience*. Elsevier.

[27] Gerald Häubl and Valerie Trifts. 2000. Consumer decision making in online shopping environments: The effects of interactive decision aids. *Marketing science* 19, 1 (2000), 4–21.

[28] Weiyin Hong, James YL Thong, and Kar Yan Tam. 2004. Designing product listing pages on e-commerce websites: an examination of presentation mode and information format. *International Journal of Human-Computer Studies* 61, 4 (2004), 481–503.

[29] Jianwei Hou and Kevin Elliott. 2021. Mobile shopping intensity: Consumer demographics and motivations. *Journal of Retailing and Consumer Services* 63 (2021), 102741.

[30] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2023. Mlagentbench: Evaluating language agents on machine learning experimentation. *arXiv preprint arXiv:2310.03302* (2023).

[31] John F Kelley. 1984. An iterative design methodology for user-friendly natural language office information applications. *ACM Transactions on Information Systems (TOIS)* 2, 1 (1984), 26–41.

[32] Anjali Khurana, Hariharan Subramonyam, and Parmit K Chilana. 2024. Why and when llm-based assistants can go wrong: Investigating the effectiveness of prompt-based interactions for software help-seeking. In *Proceedings of the 29th International Conference on Intelligent User Interfaces*. 288–303.

[33] Jesper Kjeldskov and Connor Graham. 2003. A review of mobile HCI research methods. In *International Conference on Mobile Human-Computer Interaction*. Springer, 317–335.

[34] Rafal Kocielnik, Saleema Amershi, and Paul N Bennett. 2019. Will you accept an imperfect ai? exploring designs for adjusting end-user expectations of ai systems. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–14.

[35] Seunghun Koh, Byung Hyung Kim, and Sungho Jo. 2025. Understanding the user perception and experience of interactive algorithmic recourse customization. *ACM Transactions on Computer-Human Interaction* 31, 3 (2025), 1–25.

[36] Emily Kuang, Ehsan Jahangirzadeh Soure, Mingming Fan, Jian Zhao, and Kristen Shinohara. 2023. Collaboration with conversational AI assistants for UX evaluation: Questions and how to ask them (voice vs. text). In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–15.

[37] Emily Kuang, Minghao Li, Mingming Fan, and Kristen Shinohara. 2024. Enhancing UX evaluation through collaboration with conversational AI assistants: Effects of proactive dialogue and timing. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–16.

[38] Raina Langevin, Ross J Lordon, Thi Avrahami, Benjamin R Cowan, Tad Hirsch, and Gary Hsieh. 2021. Heuristic evaluation of conversational agents. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.

[39] Jonathan Lazar, Jinjuan Heidi Feng, and Harry Hochheiser. 2017. *Research methods in human-computer interaction*. Morgan Kaufmann.

[40] Kyuwon Lee, Simone Paci, Jeongmin Park, Hye Young You, and Sylvan Zheng. [n. d.]. Applications of GPT in Political Science Research. ([n. d.]).

[41] Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. 2024. Agent Hospital: A Simulacrum of Hospital with Evolvable Medical Agents. doi:10.48550/arXiv.2405.02957 arXiv:2405.02957 [cs]

[42] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688* (2023).

[43] Xuan Liu, Jie Zhang, Haoyang Shang, Song Guo, Chengxu Yang, and Quanyan Zhu. 2024. Exploring prosocial irrationality for llm agents: A social cognition view. *arXiv preprint arXiv:2405.14744* (2024).

[44] Yuxuan Lu, Jing Huang, Yan Han, Bingsheng Yao, Sisong Bei, Jiri Gesi, Yaochen Xie, Zheshen, Wang, Qi He, and Dakuo Wang. 2025. Prompting is Not All You Need! Evaluating LLM Agent Simulation Methodologies with Real-World Online Customer Behavior Data. arXiv:2503.20749 [cs.CL] https://arxiv.org/abs/2503.20749

[45] Yuxuan Lu, Bingsheng Yao, Hansu Gu, Jing Huang, Jessie Wang, Laurence Li, Jiri Gesi, Qi He, Toby Jia-Jun Li, and Dakuo Wang. 2025. Uxagent: An llm agent-based usability testing framework for web design. *arXiv preprint arXiv:2502.12561* (2025).

[46] Ewa Luger and Abigail Sellen. 2016. "Like Having a Really Bad PA": The Gulf between User Expectation and Experience of Conversational Agents. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) *(CHI '16)*. Association for Computing Machinery, New York, NY, USA, 5286–5297. doi:10.1145/2858036.2858288

[47] Chen Luo, Dimitri Papadimitriou, Hariharan Muralidharan, Dhineshkumar Ramasubbu, Aakash Kolekar, Wenju Xu, Cong Xu, Anirudh Srinivasan, Mukesh Jain, and Qi He. 2025. Language Model Alignment for Conversational Shopping at Amazon. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Padua, Italy) *(SIGIR '25)*. Association for Computing Machinery, New York, NY, USA, 4314–4318. doi:10.1145/3726302.3731955

[48] Xanayra Marin-Lopez. 2025. Google is turning AI into a personal shopping assistant. https://www.retaildive.com/news/google-shopping-ai-personal-assistant/748893/ Accessed: 2025-06-05.

[49] Rajiv Mehta and Trishul Chilimbi. 2024. Amazon announces Rufus, a new generative AI-powered conversational shopping experience. https://www.aboutamazon.com/news/retail/amazon-rufus. Accessed: 2025-06-05.

[50] Fengran Mo, Chuan Meng, Mohammad Aliannejadi, and Jian-Yun Nie. 2025. Conversational search: From fundamentals to frontiers in the LLM era. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 4094–4097.

[51] Carol Moser, Chanda Phelan, Paul Resnick, Sarita Y Schoenebeck, and Katharina Reinecke. 2017. No such thing as too much chocolate: evidence against choice overload in e-commerce. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*. 4358–4369.

[52] San Murugesan. 2025. The rise of agentic AI: implications, concerns, and the path forward. *IEEE Intelligent Systems* 40, 2 (2025), 8–14.

[53] Isabel Briggs Myers. 2003. *MBTI manual: A guide to the development and use of the Myers-Briggs Type Indicator*. Cpp.

[54] Tahmid Nayeem and Jean Marie-IpSooching. 2022. Revisiting Sproles and Kendall's consumer styles inventory (CSI) in the 21st Century: A case of Australian consumers decision-making styles in the context of high and low-involvement purchases. *Marketing* 7, 2 (2022), 7–17.

[55] Jakob Nielsen. 1994. *Usability engineering*. Morgan Kaufmann.

[56] Nurul Zarirah Nizam and Jaafar Abdullah Jaafar. 2018. Interactive online advertising: The effectiveness of marketing strategy towards customers purchase decision. *International Journal of Human and Technology Interaction (IJHaTI)* 2, 2 (2018), 9–16.

[57] Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative Agents: Interactive Simulacra

[58] of Human Behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*. Association for Computing Machinery, New York, NY, USA, 1–22. doi:10.1145/3586183.3606763

[58] Joon Sung Park, Carolyn Q Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S Bernstein. 2024. Generative agent simulations of 1,000 people. *arXiv preprint arXiv:2411.10109* (2024).

[59] Joon Sung Park, Carolyn Q. Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S. Bernstein. 2024. Generative Agent Simulations of 1,000 People. doi:10.48550/arXiv.2411.10109 arXiv:2411.10109 [cs]

[60] Martin Porcheron, Joel E. Fischer, and Stuart Reeves. 2021. Pulling Back the Curtain on the Wizards of Oz. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW3, Article 243 (Jan. 2021), 22 pages. doi:10.1145/3432942

[61] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. ChatDev: Communicative Agents for Software Development. doi:10.48550/arXiv.2307.07924 arXiv:2307.07924 [cs]

[62] Katharina Reinecke and Abraham Bernstein. 2011. Improving performance, perceived usability, and aesthetics with culturally adaptive user interfaces. *ACM Trans. Comput.-Hum. Interact.* 18, 2, Article 8 (July 2011), 29 pages. doi:10.1145/1970378.1970382

[63] Yuqing Ren and Robert E Kraut. 2014. Agent based modeling to inform the design of multiuser systems. In *Ways of Knowing in HCI*. Springer, 395–419.

[64] Reuters. 2025. OpenAI rolls out new shopping features with ChatGPT search update. *Reuters* (28 April 2025). https://www.reuters.com/business/media-telecom/openai-rolls-out-new-shopping-features-with-chatgpt-search-update-2025-04-28/ Accessed: 2025-06-05.

[65] Lilian Rincon. 2025. Shop with AI Mode, use AI to buy and try clothes on yourself virtually. https://blog.google/products/shopping/google-shopping-ai-mode-virtual-try-on-update/. Accessed: 2025-06-05.

[66] Jennifer Rowley. 2000. Product search in e-shopping: a review and research propositions. *Journal of consumer marketing* 17, 1 (2000), 20–35.

[67] Ranjan Sapkota, Konstantinos I Roumeliotis, and Manoj Karkee. 2025. Ai agents vs. agentic ai: A conceptual taxonomy, applications and challenges. *arXiv preprint arXiv:2505.10468* (2025).

[68] Jeff Sauro and James R Lewis. 2016. *Quantifying the user experience: Practical statistics for user research*. Morgan Kaufmann.

[69] Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor. 2024. AgentClinic: A Multimodal Agent Benchmark to Evaluate AI in Simulated Clinical Environments. arXiv:2405.07960 [cs] http://arxiv.org/abs/2405.07960

[70] Eike Schneiders, Tina Seabrooke, Joshua Krook, Richard Hyde, Natalie Leesakul, Jeremie Clos, and Joel E Fischer. 2025. Objection Overruled! Lay People can Distinguish Large Language Models from Lawyers, but still Favour Advice from an LLM. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 1201, 14 pages. doi:10.1145/3706598.3713470

[71] Jing-bo Shao, Zhen-zhen Li, and Ming-ye Hu. 2014. The impact of online reviews on consumers' purchase decisions in online shopping. In *2014 International Conference on Management Science & Engineering 21th Annual Conference Proceedings*. IEEE, 287–293.

[72] Helen Sharp, Yvonne Rogers, and Jenny Preece. 2007. Interaction design: beyond human-computer interaction. *(No Title)* (2007).

[73] Yonadav Shavit, Sandhini Agarwal, Miles Brundage, Steven Adler, Cullen O'Keefe, Rosie Campbell, Teddy Lee, Pamela Mishkin, Tyna Eloundou, Alan Hickey, et al. 2023. Practices for governing agentic AI systems. *Research Paper, OpenAI* (2023).

[74] Sofia Eleni Spatharioti, David Rothschild, Daniel G Goldstein, and Jake M Hofman. 2025. Effects of LLM-based Search on Decision Making: Speed, Accuracy, and Overreliance. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–15.

[75] Samantha . Stahlke, Atiya Nova, and Pejman Mirza-Babaei. 2019. Artificial Playfulness: A Tool for Automated Agent-Based Playtesting. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems (CHI EA '19)*. Association for Computing Machinery, New York, NY, USA, 1–6. doi:10.1145/3290607.3313039

[76] Maryam Taeb, Amanda Swearngin, Eldon Schoop, Ruijia Cheng, Yue Jiang, and Jeffrey Nichols. 2024. AXNav: Replaying Accessibility Tests from Natural Language. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems (CHI '24)*. Association for Computing Machinery, New York, NY, USA, 1–16. doi:10.1145/3613904.3642777

[77] Amos Tversky and Daniel Kahneman. 1974. Judgment under Uncertainty: Heuristics and Biases: Biases in judgments reveal some heuristics of thinking under uncertainty. *science* 185, 4157 (1974), 1124–1131.

[78] Priyan Vaithilingam, Tianyi Zhang, and Elena L. Glassman. 2022. Expectation vs. Experience: Evaluating the Usability of Code Generation Tools Powered by Large Language Models. In *Extended Abstracts of the 2022 CHI Conference on*

*Human Factors in Computing Systems* (New Orleans, LA, USA) *(CHI EA '22)*. Association for Computing Machinery, New York, NY, USA, Article 332, 7 pages. doi:10.1145/3491101.3519665

[79] Dakuo Wang, Ting-Yao Hsu, Yuxuan Lu, Limeng Cui, Yaochen Xie, William Headean, Bingsheng Yao, Akash Veeragouni, Jiapeng Liu, Sreyashi Nag, et al. 2025. AgentA/B: Automated and Scalable Web A/BTesting with Interactive LLM Agents. *arXiv preprint arXiv:2504.09723* (2025).

[80] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. 2024. A Survey on Large Language Model Based Autonomous Agents. *Frontiers of Computer Science* 18, 6 (Dec. 2024), 186345. doi:10.1007/s11704-024-40231-1 arXiv:2308.11432 [cs]

[81] Lei Wang, Jingsen Zhang, Hao Yang, Zhiyuan Chen, Jiakai Tang, Zeyu Zhang, Xu Chen, Yankai Lin, Ruihua Song, Wayne Xin Zhao, et al. 2023. User behavior simulation with large language model based agents. *arXiv preprint arXiv:2306.02552* (2023).

[82] Ziyi Wang, Yuxuan Lu, Wenbo Li, Amirali Amini, Bo Sun, Yakov Bart, Weimin Lyu, Jiri Gesi, Tian Wang, Jing Huang, et al. 2025. OPeRA: A Dataset of Observation, Persona, Rationale, and Action for Evaluating LLMs on Human Online Shopping Behavior Simulation. *arXiv preprint arXiv:2506.05606* (2025).

[83] Albatool Wazzan, Stephen MacNeil, and Richard Souvenir. 2024. Comparing traditional and LLM-based search for image geolocation. In *Proceedings of the 2024 Conference on Human Information Interaction and Retrieval*. 291–302.

[84] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.

[85] Justin D. Weisz, Jessica He, Michael Muller, Gabriela Hoefer, Rachel Miles, and Werner Geyer. 2024. Design Principles for Generative AI Applications. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 378, 22 pages. doi:10.1145/3613904.3642466

[86] Michel Wermelinger. 2023. Using github copilot to solve simple programming problems. In *Proceedings of the 54th ACM Technical Symposium on Computer Science Education V. 1*. 172–178.

[87] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) *(CHI '22)*. Association for Computing Machinery, New York, NY, USA, Article 385, 22 pages. doi:10.1145/3491102.3517582

[88] Wei Xiang, Hanfei Zhu, Suqi Lou, Xinli Chen, Zhenghua Pan, Yuping Jin, Shi Chen, and Lingyun Sun. 2024. SimUser: Generating Usability Feedback by Simulating Various Users Interacting with Mobile Applications. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–17.

[89] Jun Xiao, John Stasko, and Richard Catrambone. 2007. The role of choice and customization on users' interaction with embodied conversational agents: effects on perception and performance. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 1293–1302.

[90] Chengxing Xie, Canyu Chen, Feiran Jia, Ziyu Ye, Shiyang Lai, Kai Shu, Jindong Gu, Adel Bibi, Ziniu Hu, David Jurgens, et al. 2024. Can large language model agents simulate human trust behavior? *Advances in neural information processing systems* 37 (2024), 15674–15729.

[91] Yitian Yang, Yugin Tan, Yang Chen Lin, Jung-Tai King, Zihan Liu, and Yi-Chieh Lee. 2025. Understanding How Psychological Distance Influences User Preferences in Conversational versus Web Search. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 582, 18 pages. doi:10.1145/3706598.3713770

[92] Sojeong Yun and Youn-kyung Lim. 2025. User Experience with LLM-powered Conversational Recommendation Systems: A Case of Music Recommendation. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 898, 15 pages. doi:10.1145/3706598.3713347

[93] J.D. Zamfirescu-Pereira, Richmond Y. Wong, Bjoern Hartmann, and Qian Yang. 2023. Why Johnny Can't Prompt: How Non-AI Experts Try (and Fail) to Design LLM Prompts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) *(CHI '23)*. Association for Computing Machinery, New York, NY, USA, Article 437, 21 pages. doi:10.1145/3544548.3581388

[94] Beiqi Zhang, Peng Liang, Xiyu Zhou, Aakash Ahmad, and Muhammad Waseem. 2023. Practices and challenges of using github copilot: An empirical study. *arXiv preprint arXiv:2303.08733* (2023).

[95] Yimeng Zhang, Tian Wang, Jiri Gesi, Ziyi Wang, Yuxuan Lu, Jiacheng Lin, Sinong Zhan, Vianne Gao, Ruochen Jiao, Junze Liu, et al. 2025. Shop-R1: Rewarding LLMs to Simulate Human Behavior in Online Shopping via Reinforcement Learning. *arXiv preprint arXiv:2507.17842* (2025).

[96] Hemingxi Zheng, Daibo Xiao, and Jianxin Zhou. 2025. Enhancing perceived value in human-computer interaction: The mediating role of user participation and the moderating role of task complexity. *International Journal of Human–Computer Interaction* 41, 7 (2025), 4261–4270.

[97] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Yonghao Wu, Shixiang Zhuang, Zi Lin, Zhen Li, Dacheng Li, Eric Xing, Ion Stoica, and Joseph E. Gonzalez. 2023. MT-Bench: A Multi-Turn Benchmark for Evaluating LLMs as Chatbots. In *Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track.* https://arxiv.org/abs/2306.05685

[98] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems* 36 (2023), 46595–46623.

[99] Qingxiao Zheng, Minrui Chen, Pranav Sharma, Yiliu Tang, Mehul Oswal, Yiren Liu, and Yun Huang. 2025. EvAlignUX: Advancing UX Evaluation through LLM-Supported Metrics Exploration. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 1051, 25 pages. doi:10.1145/3706598.3714045

[100] Mingyuan Zhong, Ruolin Chen, Xia Chen, James Fogarty, and Jacob O Wobbrock. 2025. ScreenAudit: Detecting Screen Reader Accessibility Errors in Mobile Apps Using Large Language Models. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–19.

[101] Shenzhe Zhu, Jiao Sun, Yi Nian, Tobin South, Alex Pentland, and Jiaxin Pei. 2025. The Automated but Risky Game: Modeling Agent-to-Agent Negotiations and Transactions in Consumer Markets. arXiv:2506.00073 [cs.AI] https://arxiv.org/abs/2506.00073

[102] Mingchen Zhuge, Changsheng Zhao, Dylan Ashley, Wenyi Wang, Dmitrii Khizbullin, Yunyang Xiong, Zechun Liu, Ernie Chang, Raghuraman Krishnamoorthi, Yuandong Tian, et al. 2024. Agent-as-a-judge: Evaluate agents with agents. *arXiv preprint arXiv:2410.10934* (2024).

# A  Surveys

## A.1  User Persona Survey

### A.1.1  Part 1: Demographic Information.

(1) Please state your gender: Male, Female, Non-binary or gender non-conforming, Prefer to self-describe, Prefer not to say

(2) What is your age? Under 18, 18–24, 25–34, 35–44, 45–54, 55–64, 65+

(3) Which city and state do you live in? (e.g., San Diego, California)

(4) What is the highest level of education you have completed? High school or less, High school diploma or GED, Some college (no degree), Associate or technical degree, Bachelor's degree, Graduate or professional degree, Prefer not to say

(5) What was your total household income before taxes during the past 12 months? Less than $25,000, $25,000–$49,999, $50,000–$74,999, $75,000–$99,999, $100,000–$149,999, $150,000 or more, Prefer not to say

(6) What best describes your employment status over the last three months? Full-time employee, Part-time employee, Self-employed, Unemployed and looking for work, Student, Retired, Other

(7) Do you live alone or live with others? If so, who are they? (Optional)

(8) Use three sentences to describe yourself. (Free text)

(9) Use three sentences to describe your daily routine. (Free text)

### A.1.2  Part 2: Shopping Habits.

(1) How often do you shop online? More than three times a week, Once to twice a week, Once every couple of weeks, Less than once a month

(2) How much money (in USD) do you spend on online shopping per month (not including food or delivery services)?

(3) Do you have a paid membership for expedited delivery (e.g., Prime)? Yes / No

Confidence in product selection (1 = Not at all confident, 5 = Extremely confident):

(1) Monitor (e.g., technical specs)
(2) Chair for your home (e.g., ergonomics, brands)
(3) Summer outfit (e.g., style, brands)
(4) Jacket (e.g., durability, weather protection)

Shopping attitudes (1 = Strongly Disagree, 5 = Strongly Agree):

(1) I tend to shop more during holidays (e.g., Black Friday, holiday sales).
(2) Online ads attract my attention and are a good source of information.
(3) I usually do a lot of research (e.g., reading reviews) before making a purchase.
(4) I prioritize delivery speed and delivery fee of the product.
(5) Getting high-quality online products is very important for me.
(6) The more expensive brands are usually my choice.
(7) The more I learn about online products, the harder it seems to choose the best.
(8) I shop quickly, buying the first product or brand that seems good enough.
(9) Once I find a brand I like, I stick with it.
(10) I would buy a new or different brand just to see what it is like.
(11) I enjoy shopping for online products just for the fun of it.
(12) I look carefully to find the best value for money when shopping online.

### A.1.3 Part 3: Personality Test.

*Instruction.* Please read each statement and indicate how well it describes you. Use a 5-point scale: 1 = Very Inaccurate, 2 = Moderately Inaccurate, 3 = Neither, 4 = Moderately Accurate, 5 = Very Accurate. For each item below, respondents select one value from 1–5.

(1) I am the life of the party
(2) Feel little concern for others
(3) Am always prepared
(4) Get stressed out easily
(5) Have a rich vocabulary
(6) Don't talk a lot
(7) Am interested in people
(8) Leave my belongings around
(9) Am relaxed most of the time
(10) Have difficulty understanding abstract ideas
(11) Feel comfortable around people
(12) Insult people
(13) Pay attention to details
(14) Worry about things
(15) Have a vivid imagination
(16) Keep in the background
(17) Sympathize with others' feelings
(18) Make a mess of things
(19) Seldom feel blue
(20) Am not interested in abstract ideas
(21) Start conversations
(22) Am not interested in other people's problems
(23) Get chores done right away
(24) Am easily disturbed
(25) Have excellent ideas
(26) Have little to say
(27) Have a soft heart
(28) Often forget to put things back in their proper place
(29) Get upset easily
(30) Do not have a good imagination
(31) Talk to a lot of different people at parties
(32) Am not really interested in others
(33) Like order
(34) Change my mood a lot
(35) Am quick to understand things
(36) Don't like to draw attention to myself
(37) Take time out for others
(38) Shirk my duties
(39) Have frequent mood swings
(40) Use difficult words
(41) Don't mind being the center of attention
(42) Feel others' emotions
(43) Follow a schedule
(44) Get irritated easily
(45) Spend time reflecting on things
(46) Am quiet around strangers
(47) Make people feel at ease
(48) Am exacting in my work
(49) Often feel blue
(50) Am full of ideas

MBTI (Optional): What is your MBTI personality type? (Free-text)

## A.2 Post-Study Survey

### A.2.1 General Satisfaction.

(1) Did you chat with the AI shopping assistant during the shopping task?
*Response:* Yes / No

(2) Overall satisfaction with the final product you picked:
   1 = Very Unsatisfied, 2 = Somewhat Unsatisfied, 3 = Neutral,
   4 = Somewhat Satisfied, 5 = Very Satisfied

*A.2.2  Interaction Experience.  Instruction.* Please indicate your agreement with the following statements (1 = Strongly Disagree, 5 = Strongly Agree).

   (1) The assistant's responses matched my query.
   (2) The conversation with the assistant felt engaging.
   (3) The assistant's responses were logical and coherent.
   (4) The assistant's responses matched product information.
   (5) It was easy to interact with the assistant via chat.
   (6) I found the interaction enjoyable.
   (7) The assistant was helpful for me to buy the product.
   (8) I would love to use the assistant to shop in the future.
   (9) I will recommend others to use the assistant to shop.
   (10) I found the interaction mentally demanding.
   (11) I prefer the assistant over traditional search or manual browsing.

*A.2.3  Trust and Concerns.  Instruction.* Please indicate your agreement with the following statements (1 = Strongly Disagree, 5 = Strongly Agree).

   (1) I trust the responses provided by the assistant.
   (2) I was concerned that the assistant may manipulate me using sponsored recommendations.

*A.2.4  Open-ended Questions.*

   (1) Can you describe a time when the assistant was particularly helpful or frustrating during your shopping experience?
   (2) In the future, as AI shopping assistants become more common, will you change your online shopping behavior and how?

Table 4: Participant Demographics

| PID | Age | Sex | Ethnicity | Country of Residence | Language | Student Status | Employment Status |
|---|---|---|---|---|---|---|---|
| P1 | 33 | Female | White | United States | English | No | Part-Time |
| P2 | 56 | Male | White | United States | English | Unknown | Unknown |
| P3 | 21 | Female | White | United States | English | Yes | Part-Time |
| P4 | 56 | Female | White | United States | English | Unknown | Unknown |
| P5 | 29 | Male | White | United States | English | No | Unemployed (and job seeking) |
| P6 | 48 | Female | White | United States | English | Unknown | Unknown |
| P7 | 31 | Female | White | United States | English | Unknown | Unknown |
| P8 | 41 | Female | White | United States | English | No | Unemployed (and job seeking) |
| P9 | 37 | Male | White | United States | English | No | Part-Time |
| P10 | 32 | Male | White | United States | English | Unknown | Full-Time |
| P11 | 52 | Male | White | United States | English | Unknown | Unknown |
| P12 | 54 | Female | Black | United States | English | Unknown | Unknown |
| P13 | 36 | Male | Black | United States | English | Unknown | Unknown |
| P14 | 32 | Female | Black | United States | English | Yes | Full-Time |
| P15 | 56 | Male | Mixed | United States | English | Yes | Part-Time |
| P16 | 41 | Female | White | United States | English | No | Full-Time |
| P17 | 49 | Male | White | United States | English | Unknown | Unknown |
| P18 | 35 | Female | White | United States | English | No | Full-Time |
| P19 | 30 | Male | White | United States | English | Unknown | Unknown |
| P20 | 27 | Female | Mixed | United States | English | Unknown | Unknown |
| P21 | 50 | Male | White | United States | English | No | Unknown |
| P22 | 37 | Female | Black | United States | English | No | Full-Time |
| P23 | 38 | Male | Black | United States | English | Unknown | Unknown |
| P24 | 49 | Male | Asian | United States | Chinese | No | Full-Time |
| P25 | 30 | Female | White | United States | English | No | Other |
| P26 | 55 | Male | White | United States | English | No | Full-Time |
| P27 | 50 | Male | White | United States | English | Unknown | Unknown |
| P28 | 39 | Male | Asian | United States | English | No | Part-Time |
| P29 | 59 | Female | Black | United States | English | No | Other |
| P30 | 30 | Female | White | United States | English | No | Part-Time |
| P31 | 26 | Male | White | United States | English | Unknown | Other |
| P32 | 24 | Female | Black | United States | English | Yes | Full-Time |
| P33 | 34 | Female | White | United States | English | Unknown | Full-Time |
| P34 | 33 | Female | White | United States | English | No | Unknown |
| P35 | 35 | Female | White | United States | English | No | Unemployed (and job seeking) |

Table 4: Participant Demographics (Continued)

| PID | Age | Sex | Ethnicity | Country of Residence | Language | Student Status | Employment Status |
|-----|-----|-----|-----------|---------------------|----------|----------------|-------------------|
| P36 | 39 | Male | Other | United States | English | No | Full-Time |
| P37 | 49 | Female | White | United States | English | No | Full-Time |
| P38 | 42 | Male | White | United States | English | No | Unemployed (and job seeking) |
| P39 | 32 | Female | White | United States | English | No | Unknown |
| P40 | 31 | Female | White | United States | English | No | Unemployed (and job seeking) |